

Linked open data og geografiske autoritetsdata

Knut Anton Bøckman
Stig Svenningsen

Det Kongelige Bibliotek
HisKIS-seminar SDU 2016-09-14

Match af
lokale
registreringer med
publicerede
data:

- geonames
- open street maps
- GSTs AWS service
- Stednanvebase

Identifikation

KBB01-sys	621*a	621-a-underinddeling	url_OSM
000018748	København	Nyhavn	http://www.openstreetmap.org/node/894229676
000019987	København	Nørrebrogade, Nørrebro Runddel	http://www.openstreetmap.org/node/664257007
000019985	København	Nørrebrogade, Nørrebro Runddel	http://www.openstreetmap.org/node/664257007
000019976	København	Nørrebrogade	http://www.openstreetmap.org/way/310869916
000019963	København	Nørrebrogade	http://www.openstreetmap.org/way/310869916
000019950	København	Vester Voldgade	http://www.openstreetmap.org/way/149756673
000019918	København	Frederiksborggade edit	http://www.openstreetmap.org/way/77628809
000019856	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019853	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019853	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000019846	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019846	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000019827	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019804	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019804	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000019804	København	Tomebuskegade	http://www.openstreetmap.org/way/1931867
000019698	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000019683	København	Nørre Voldgade, Nørreport Station	http://www.openstreetmap.org/node/2343875251

Agenda: forbinde datasæt

- Formål – bedre dataudnyttelse
- Metode – Linked Open Data
- Værktøj – OpenRefine
- Kilder – geografiske autoritetsdata
- Udfordringer – skala, granularitet

Bibliotek

Indsamle Bevare Formidle

digitalisering

≠

formidling



Use case:
Jeg ønsker at finde et historisk
billede fra slagterierne i Kødbyen.

Use case

Jeg ønsker at finde et historisk billede fra slagterierne i Kødbyen.

- Mon ikke Det Kongelige Bibliotek har det?
- Jeg finder Det Kongelige Biblioteks hjemmeside
- Jeg finder søgesystemet
- Jeg søger på Halmtorvet

Alle

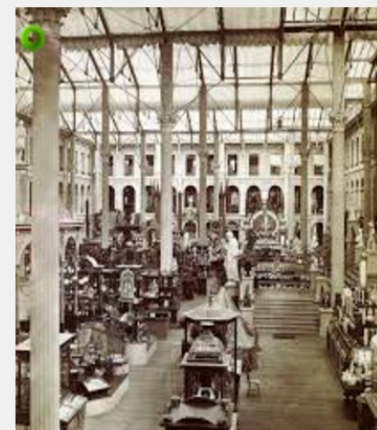
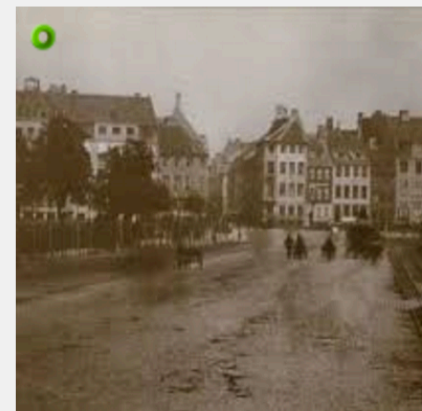
Billeder

Maps

Videoer

Mere ▾

Søgeværktøjer



Digital formidling

- Materiale: Digitalt og publiceret
- Tilgængelig for indeksering
- Kan findes via almene værktøj
- Alment anvendte beskrivelser
- Tilpas formidling til indhold

nogle eksempler

DPLA

- <https://dp.la/map?>

Haverford

- <http://cope.haverford.edu/letterNetwork>

Icelandic Saga Map

- <http://sagamap.hi.is/is/>

NYPL

- <http://pic.nypl.org>

Biblioteksdata transformeres til
MERE ANVENDELIGE DATA

Tim Berners-Lee: Linked data (2006)

1. Use URIs as names for **things**
2. Use **HTTP URIs** so that people can look up those names.
3. When someone looks up a URI, **provide useful information**, using the **standards** (RDF*, SPARQL)
4. Include **links** to other URIs, so that they can **discover more things**.

Linkning - http

Web of Documents (html)

- Links mellem dokumenter
- Information i tekst-streng
- Semantik ligger i kontekst
- Lavet for mennesker; ikke maskinevenlig

Web of Data (RDF)

- Entydig reference for ting
- Entydig reference for relation
- Brug af anerkendte internet-standarder

Data

- Resource Description Framework (RDF)
- Data beskrives i prædikative udsagn (triples)
 - Subjekt (Moesgård museum)
 - relation (ligger i)
 - Objekt (Århus)
- Selvstændige, meningsfulde **enheder**
 - **Uafhængig** af dokumentkontekst
 - Forøget betydning ved **linkning** 😊

Strings to things



- Erstat tekstbaserede navne (subjekt +objekt) med URI'er
- Også URI for prædikatet (relationen mellem dem)

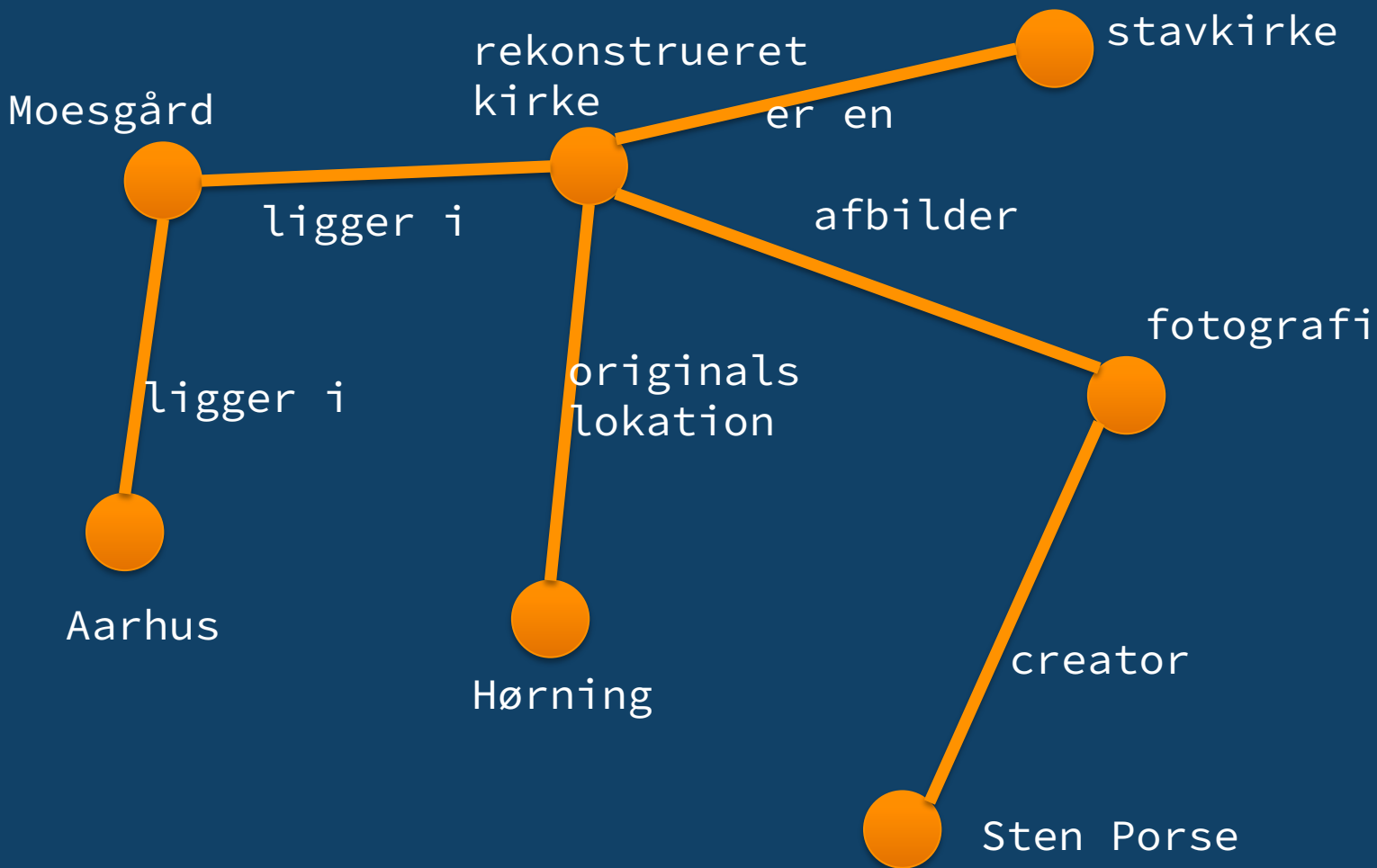
Moesgård Museum er placeret i Aarhus

`<https://www.wikidata.org/wiki/
Q3299384>`

`<http://purl.org/dc/terms/Location>`

`<http://vocab.getty.edu/tgn/7003442>`

Linkning (tænkt eks.)



Bibliotekets beskrivelse

- Dokument der beskriver en ressource
- Format: Danmarc2
- Kodede indførsler
- Semantisk scope er dokumentet
- Data som tekststreng

	<u>l</u>	skan TÜRCK 50344
<u>248</u>	<u>00</u>	<u>g</u> DT103789
	<u>a</u>	Kølehus. Interiør, svin
	<u>e</u>	foto. Sven TÜRCK
	<u>l</u>	neg. TÜRCK 43223
<u>248</u>	<u>00</u>	<u>g</u> DT103790
	<u>a</u>	Kølehus. Interiør
	<u>e</u>	foto. 152470
<u>300</u>	<u>00</u>	<u>n</u> Fotografi
<u>599</u>	<u>00</u>	<u>b</u> KBB980603
<u>621</u>	<u>00</u>	<u>a</u> København, Halmtorvet
	<u>b</u>	Danmark
<u>700</u>	<u>00</u>	<u>a</u> TÜRCK
	<u>h</u>	Sven
	<u>f</u>	fotograf
	<u>c</u>	ca. 1897-1954
<u>700</u>	<u>00</u>	<u>a</u> Nielsen

Biblioteksdata (billeder)

- Strukturerede data
- Indeholder stedsbeskrivelser
- Mangler linkning til andre data
- Ikke alment udbredt format
- Ikke entydige navne på ting

Ting - navn: URI

- Eksport af kilddedata
 - To datasæt: København og udenfor
 - CSV
- Normalisering/tilpasning
 - OpenRefine
- Match mod autoritetskilder
 - Geonames, OpenStreetMaps, andre
- Genimport af berigede data

Håndtering af data

Eksport: csv

Load i OpenRefine

Et google regneark
med (meget) ekstra
funktionalitet

Grafisk UI: lav
tærskel

Vedligeholdes af
aktiv community

000001035		aSlangerup	bDanmark
000001350		aMünchen	bTyskland
000002435		aHälsingborg	bSverige
000002635		aSchloss Husum	bTyskland
000002735		aHobro	bDanmark
000002835		aFrydenlund Slot	bDanmark
000003935		aSønderborg	bDanmark
000003935		aTinglev	bDanmark
000004735		aKarise, Karise kirke	bDanmark
000006235			bDanmark
000006350		aWashington	bUSA

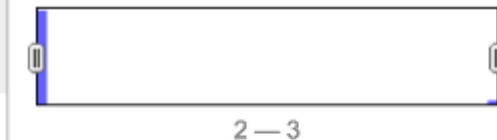
Normalisering af data

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

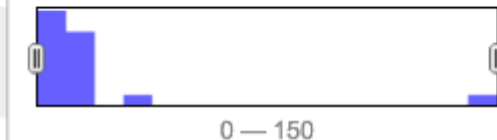
Method key collision Keying Function ngram-fingerprint Ngram Size 4 18 clusters found

3	29	<ul style="list-style-type: none">H. C. Andersens Boulevard (27 rows)H.C. Andersens Boulevard (1 rows)H.C.Andersens Boulevard (1 rows)	<input type="checkbox"/>	H. C. Andersens Boulevard
2	6	<ul style="list-style-type: none">Gammel Kalkbrænderi Vej (4 rows)Gammel Kalkbrænderivej (2 rows)	<input type="checkbox"/>	Gammel Kalkbrænderi Vej
2	10	<ul style="list-style-type: none">Kristen Bernikows Gade (5 rows)Kristen Bernikowsgade (5 rows)	<input type="checkbox"/>	Kristen Bernikows Gade
2	4	<ul style="list-style-type: none">Lille Helliggeist Stræde (3 rows)Lille Helliggeiststræde (1 rows)	<input type="checkbox"/>	Lille Helliggeist Stræde
2	7	<ul style="list-style-type: none">Kastellet, Kastelskirken (6 rows)Kastellet,Kastelskirken (1 rows)	<input type="checkbox"/>	Kastellet, Kastelskirken
2	4	<ul style="list-style-type: none">Amagertorv, Hellig Geistes Kirke (3 rows)Amagertorv, Helliggeistes Kirke (1 rows)	<input type="checkbox"/>	Amagertorv, Hellig Geistes Kir
2	2	<ul style="list-style-type: none">H. P. Ørums Gade (1 rows)H. P. Ørumsgade (1 rows)	<input type="checkbox"/>	H. P. Ørums Gade
2	8	<ul style="list-style-type: none">Nørrevold (7 rows)Nørrevold (1 rows)	<input type="checkbox"/>	Nørrevold

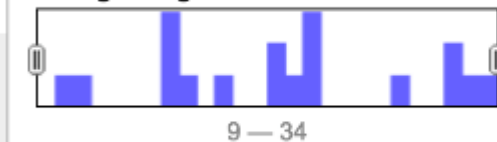
Choices in Cluster



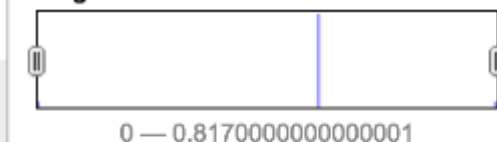
Rows in Cluster



Average Length of Choices



Length Variance of Choices



Normalisering af data

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

Method

Keying Function

191 clusters found

- [Frederiksholms Kanal, Fæstningsmaterialgården](#) (1 rows)
- [Frederiksholms Kanal, Hestegårdens Kaserne](#) (1 rows)
- [Frederiksholms Kanal, Materialgården](#) (1 rows)
- [Frederiksholms Kanal, Nationalmuseets Klunkehjem](#) (1 rows)
- [Frederiksholms Kanal, Wedells Palæ](#) (1 rows)

16

207

- [Kongens Nytorv](#) (64 rows)
- [Kongens Nytorv, Det kongelige Teater](#) (47 rows)
- [Kongens Nytorv, Charlottenborg](#) (31 rows)
- [Kongens Nytorv, Hotel d'Angleterre](#) (13 rows)
- [Kongens Nytorv, Hovedvagten](#) (10 rows)
- [Kongens Nytorv, Magasin du Nord](#) (9 rows)
- [Kongens Nytorv, Krinsen](#) (8 rows)
- [Kongens Nytorv, Gjethuset](#) (5 rows)
- [Kongens Nytorv, Thotts Palæ](#) (5 rows)
- [Kongens Nytorv, Harsdorffs Palæ](#) (4 rows)
- [Kongens Nytorv, Kanneworffs Hus](#) (3 rows)
- [Kongens Nytorv, Stephan à Porta](#) (3 rows)
- [Kongens Nytorv, Store Nordiske Telegrafsekskab](#) (2 rows)
- [Kongens Nytorv, Charlottenborg, Udstillingsbygningen](#) (1 rows)
- [Kongens Nytorv, Den Militær Højskole](#) (1 rows)
- [Kongens Nytorv, Å Porta](#) (1 rows)



Kongens Nytorv

11

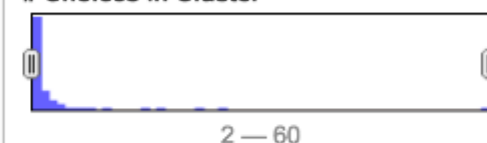
346

- [Vesterbrogade, Tivoli](#) (177 rows)
- [Vesterbrogade](#) (151 rows)
- [Vesterbrogade, Industriforeningen](#) (9 rows)
- [Vesterbrogade, Fiksdalstøtten](#) (9 rows)

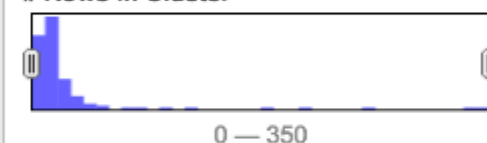


Vesterbrogade, Tivoli

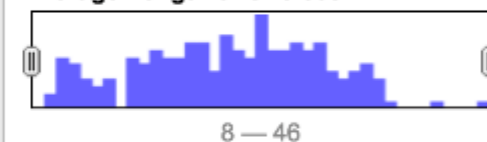
Choices in Cluster



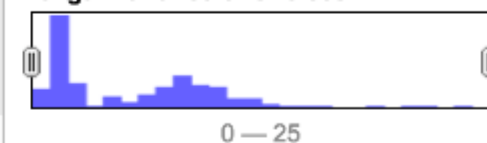
Rows in Cluster




Average Length of Choices





Length Variance of Choices





Match data mod andre kilder


Freebase Query-based Reconciliation 


NamedEntity 


Sindice 

LCSH (backup) 

Virtual International Authority File 

VIAF (via refine.codefork.com) 

GeoNames Reconciliation Service 

CSV Reconciliation service 

» Access [Service API](#)

Reconcile each cell to an entity of one of these types:

- ☒ /geonames/name
- ☐ /geonames/name_startsWith
- ☐ /geonames/name_equals
- ☐ /geonames/all

Also use relevant details from other columns:

Column	Include?	As Property
sys	<input type="checkbox"/>	<input type="text"/>
felt	<input type="checkbox"/>	<input type="text"/>
I	<input type="checkbox"/>	<input type="text"/>
\$by	<input type="checkbox"/>	<input type="text"/>
\$land	<input type="checkbox"/>	<input type="text"/>
\$område	<input type="checkbox"/>	<input type="text"/>
\$andet	<input type="checkbox"/>	<input type="text"/>
Column	<input type="checkbox"/>	<input type="text"/>
fetchet data	<input type="checkbox"/>	<input type="text"/>
land	<input type="checkbox"/>	<input type="text"/>
område	<input type="checkbox"/>	<input type="text"/>

- ☐ Reconcile against type:
- ☐ Reconcile against no particular type
- ☒ Auto-match candidates with high confidence

Hent data ind

OpenStreetMaps

GeoNames

						▼ KBB01-sys	▼ 621*a	▼ 621-a-underinddeling	▼ url_OSM
						000018748	København	Nyhavn	http://www.openstreetmap.org/node/894229676
						000019987	København	Nørrebrogade, Nørrebros Runddel	http://www.openstreetmap.org/node/664257007
						000019985	København	Nørrebrogade, Nørrebros Runddel	http://www.openstreetmap.org/node/664257007
						000019976	København	Nørrebrogade	http://www.openstreetmap.org/way/310869916
000000383						000019963	København	Nørrebrogade	http://www.openstreetmap.org/way/310869916
						000019950	København	Vester Voldgade	http://www.openstreetmap.org/way/149756673
000000384						000019918	København	Frederiksborggade edit	http://www.openstreetmap.org/way/77628809
						000019856	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000000878						000019853	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
						000019853	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000001033						000019846	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
						000019846	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000001034						000019827	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000001119						000019804	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
						000019804	København	Fiolstræde	http://www.openstreetmap.org/way/12645065
000001176						000019804	København	Tornebuskegade	http://www.openstreetmap.org/way/1931867
						000019698	København	Nørre Voldgade	http://www.openstreetmap.org/way/225625919
000001177						000019683	København	Nørre Voldgade, Nørreport Station	http://www.openstreetmap.org/node/2343875251
000001195						Bulbjerg 57.1575, 9.02588 http://sws.geonames.org/2623248			Danmark
000001208						Frederiksborg Slot 55.93451, 12.30041 http://sws.geonames.org/2621938			Danmark

Kilder til geografiske

AUTORITETSDATA

Kilder (1/2): LOD

Geonames

- Åben API
- Nem rekonciliering
- Åbne data
- Linked data
- Stabil URI
- Dårlig detaljegrad for gadenavne i København

Open Street Maps

- Begrænset API
- Åbne data
- Leverer data i JSON, skal parses
- Ingen URI-policy
- God detaljegrad for gadenavne i København
- Variabelt datagrundlag

Kilder (2/2)

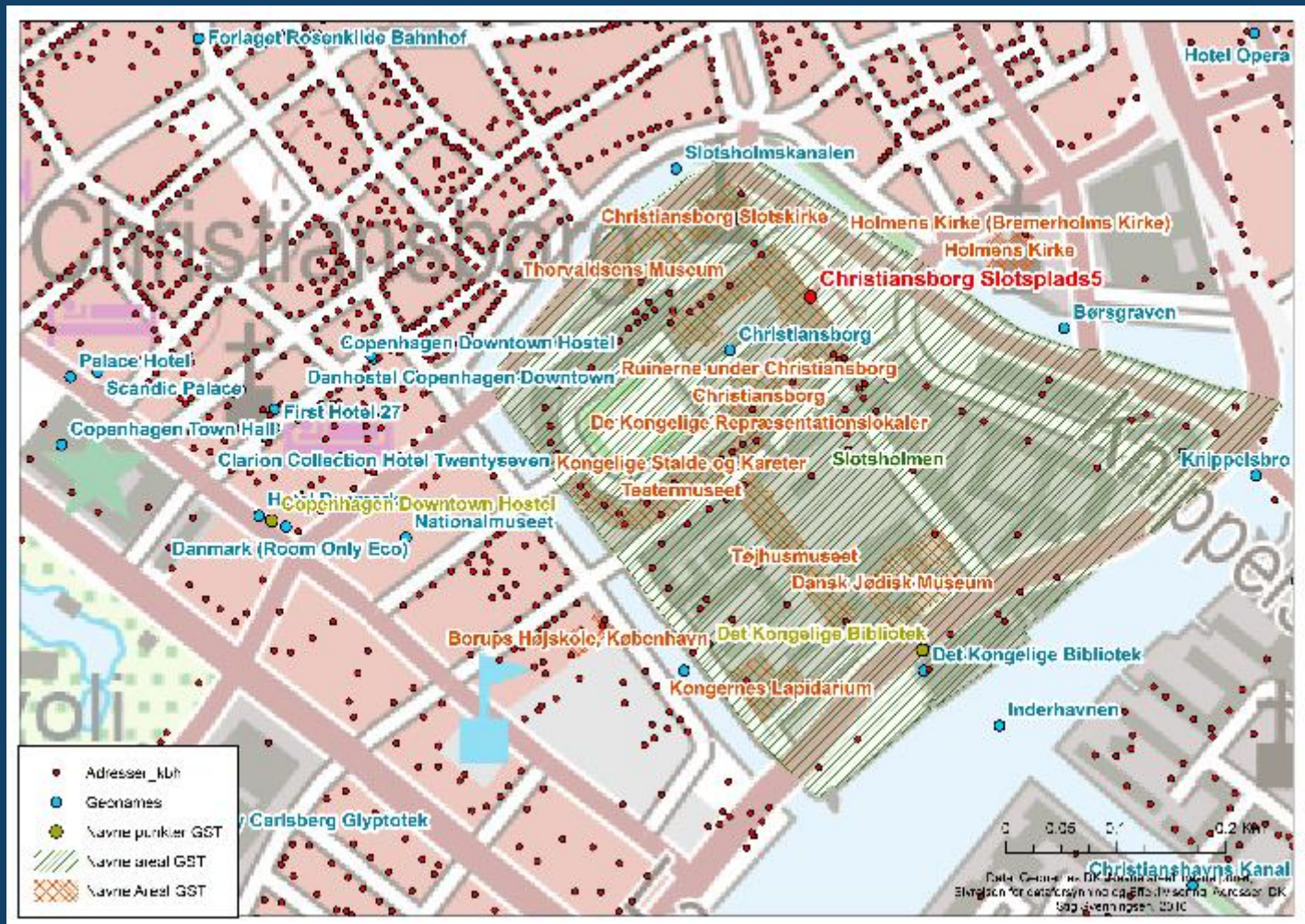
GST adresse-udtræk
fra aws.dk

- Ikke linked data
- Ingen åben API
- Data åbne?
- Datadump med identificerede indførsler på husnummer-niveau

Stednavnebase fra
Afdeling for
navneforskning

- Historiske stednavne
- Meget detaljeret udenfor København
- Ikke detaljeret i København
- Ikke publiceret
- Ikke åbne data
- Ikke linked data

Udfordringer med LOD - fra en geografisk perspektiv



Udfordringer med brugen af geografisk sted data

- Geometri for den geografiske reference
 - Punkt.
 - Linje.
 - Polygon.
- Skala
 - Christiansborg: område, stednavn, begreb, adresse(r)
- Præcision
 - Det Kongelige Bibliotek i to punkter.
- Bias
 - Geonames mest butikker.

referencer

- Tim Berners-Lee (2006): Linked data
<https://www.w3.org/DesignIssues/LinkedData.html>
- OpenRefine: <http://openrefine.org>
- Blog om LOD i DEFF-projektet Linked, Open & Social: <http://los3blog.wordpress.com/>
- Se også eksemplerne i slide 9.